# 5 minute Summary Lecture - 1

### Introduction to Unsupervised Learning

Unsupervised learning is a type of machine learning where the algorithm learns patterns from untagged data. The goal is to model the underlying structure or distribution in the data to learn more about it. Unlike supervised learning, unsupervised learning doesn't work with predictions or labels; instead, it focuses on finding relationships within the data.

### Clustering and K-Means Algorithm

Clustering is a significant part of unsupervised learning aimed at grouping sets of objects in such a way that objects in the same group (called a cluster) are more similar to each other than to those in other groups. K-Means is one of the simplest and most popular clustering algorithms. It aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean.

### The Math Behind K-Means

K-Means clustering works by initializing k centroids randomly, then iterating over two steps:

1. **Assignment step**: Assign each data point to the nearest centroid.

2. **Update step**: Update the centroids to be the mean of the points assigned to them.

The process repeats until the assignments no longer change or the changes are below a certain threshold, indicating that the algorithm has converged.

## Lloyd's Algorithm

Lloyd's Algorithm, often synonymous with K-Means, refers to the two-step iterative approach K-Means uses to converge on a solution. It's worth noting that while efficient, Lloyd's algorithm can sometimes fall into local minima, and its performance can be sensitive to the initial choice of centroids.

## Implementing K-Means from Scratch

When implementing K-Means from scratch, the key steps involve initializing centroids randomly, then iteratively updating the cluster assignments and the centroids until the algorithm converges. This process involves calculating distances between data points and centroids, typically using Euclidean distance, and recalculating cluster centroids after every assignment step.

## Determining the Optimal Number of Clusters (K)

Determining the right number of clusters, K, is crucial for K-Means performance. Techniques like the Elbow Method involve plotting the within-cluster sum of squares (WCSS) against the number of clusters and looking for the "elbow" point where the rate of decrease sharply changes. The Silhouette Score is another method, provides a measure of how similar an object is to its own cluster compared to other clusters.

## Practical Application: Customer Segmentation

Applying K-Means for customer segmentation involves preprocessing steps like feature scaling, followed by fitting the K-Means model to the data and interpreting the resulting clusters. Each cluster can represent a different customer segment, which can then be targeted with tailored marketing strategies.

## Evaluation and Challenges

Evaluating unsupervised learning models like K-Means can be challenging due to the absence of ground truth labels. Metrics like the Dunn Index, which involves the ratio of the smallest distance between observations not in the same cluster to the largest intra-cluster distance, can be used. However, the ultimate measure often involves domain-specific knowledge and the practical usefulness of the clustering results.